

# VQ-TEGAN: Data Augmentation with Text Embeddings for Few-shot Learning

Anonymous ACL submission

## Abstract

Data augmentation is crucial for the fine-tuning of pre-trained models and the optimization of limited data utilization, particularly within the realm of few-shot learning. Traditionally, these techniques have been applied at the word and sentence levels, with little research conducted within the embedding space. This paper introduces **VQ-TEGAN**, a novel data augmentation approach designed to generate embeddings specifically for a few-shot learning. VQ-TEGAN generates embeddings that augment the few-shot dataset by training directly within the PLMs' word embedding, employing a customized loss function. Empirical validation on GLUE benchmark datasets demonstrates that VQ-TEGAN markedly improves text classification performance. Additionally, we investigate the application of VQ-TEGAN with RoBERTa-large and BERT-large, offering insight for further application.

## 1 Introduction

Text classification is a crucial task in natural language processing (NLP) (Kowsari et al., 2019). Although fine-tuning pre-trained language models (PLMs) on large datasets is highly effective, performance declines with smaller training data sizes (Gao et al., 2020; Longpre et al., 2020). This is due to the lack of diverse examples. Data augmentation has emerged as a solution to improve model performance with limited data, applicable in various fields such as healthcare (Eaton-Rosen et al., 2018; Ker et al., 2017), finance (Fons et al., 2020; El-Laham and Vyetrenko, 2022), and computer vision (Zhang et al., 2017; Chen et al., 2020b).

In NLP, data augmentation is often performed through word-level manipulation (*e.g.*, EDA (Wei and Zou, 2019) and AEDA (Karimi et al., 2021)). Recent advances include sentence-level interpolation methods like MixText (Zhang et al., 2022) and Treemix (Zhang et al., 2022; Chen et al., 2020a).

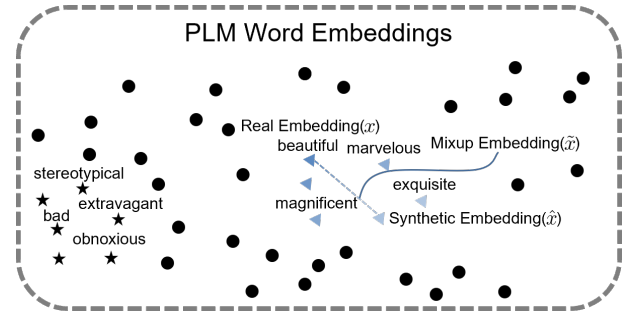


Figure 1: Graphical abstract of VQ-TEGAN. The primary aim of VQ-TEGAN is to produce synthetic embeddings that closely approximate the original real embeddings. Subsequently, the synthetic embedding is mixed with the real embedding to formulate a mixup embedding, which resides within a space comparable to that of other synonymous embeddings.

In addition, language-model-based augmentations such as LAMBADA (Anaby-Tavor et al., 2020), BF-Translation (Body et al., 2021), BART ProtAugment (Dopierre et al., 2021), and SSMBA (Ng et al., 2020) have been developed. While LAMBADA and BART ProtAugment require separate fine-tuning for data augmentation, SSMBA and BF-Translation do not, but they demand significant storage space and time due to the need for large language models or the Google Translation API.

Before training a language model, sentences are tokenized and converted to embeddings, which are used as direct input (Mikolov et al., 2013). Some works have applied data augmentation at the embedding level. For example, Wang and Yang (2015) used semantic and lexical embeddings from Word2Vec (Mikolov et al., 2013) to replace original words with  $k$ -nearest neighbor vectors. TreeMixup (Guo et al., 2019) applies linear interpolation to word and sentence embeddings, pioneering this technique in NLP tasks. TACLR (Jia et al., 2023) combines TreeMixup and EDA for contrastive learning. Recent studies show promising results using models that generate synthetic sentence

065 embeddings similar to real sentences (Onan, 2023;  
066 Jian et al., 2022). These methods effectively en-  
067 hance text embeddings to supplement insufficient  
068 data.

069 This research proposes Vector-Quantized Text  
070 Embedding Generative Adversarial Networks (VQ-  
071 TEGAN). VQ-TEGAN leverages the capabilities  
072 of Vector Quantized Generative Adversarial Net-  
073 work (VQ-GAN) (Esser et al., 2021) to generate  
074 text embeddings optimized for the semantic repre-  
075 sentation provided by word embeddings in PLMs  
076 (e.g., RoBERTa-large (Liu et al., 2019) and BERT-  
077 large (Devlin et al., 2018)). VQ-TEGAN is based  
078 on the understanding that the word embeddings  
079 of PLMs can capture deep linguistic properties  
080 beyond simple syntactic structures. We hypoth-  
081 esize that synthetic embeddings generated by VQ-  
082 TEGAN can encapsulate complex features such as  
083 context and sentiment, crucial for few-shot learn-  
084 ing tasks. Synthetic embeddings are employed in  
085 PLM training to provide new text examples that pre-  
086 serve semantic consistency and syntactic accuracy  
087 with the few-shot embedding data. This approach  
088 aligns with Brown et al. (2020), demonstrating that  
089 language models trained on extensive datasets can  
090 leverage prior knowledge to perform tasks with  
091 limited examples.

092 Our contributions can be summarized as follows:

- 093 • We propose a novel data augmentation model,  
094 VQ-TEGAN, for generating synthetic embed-  
095 dings located in a similar space as real embed-  
096 dings as illustrated in Figure 1.
- 097 • VQ-TEGAN is a lightweight model for data  
098 augmentation, allowing easy application and  
099 minimal storage requirements.
- 100 • We introduce a novel loss function suitable for  
101 NLP embeddings to train VQ-TEGAN.
- 102 • Experimental results indicate that VQ-  
103 TEGAN outperforms benchmarks in few-shot  
104 learning.
- 105 • The adequacy of the generated embeddings is  
106 confirmed by analyzing their meaning using  
107 cosine similarity to the word embeddings in  
108 PLMs.

## 109 2 Related Work

### 110 2.1 Generative Model

111 The evolution of generative models has been led  
112 by the advances of autoencoders (Ranzato et al.,  
113 2007). Variational Autoencoders (VAE) (Kingma  
114 and Welling, 2013) use neural networks to en-

115 code input data into a lower-dimensional latent  
116 space and decode it back, optimizing the lower  
117 bound on the likelihood of the data. This enables  
118 tasks such as data generation and feature extraction.  
119 Generative Adversarial Networks (GAN) (Good-  
120 fellow et al., 2014) employ two neural networks,  
121 a generator and a discriminator, training them si-  
122 multaneously in a competitive setting to generate  
123 data samples that are indistinguishable from real  
124 data. Wasserstein GAN (WGAN) (Arjovsky et al.,  
125 2017) improves on traditional GANs by using a  
126 Wasserstein distance metric for the loss function,  
127 improving training stability and addressing mode  
128 collapse, resulting in higher-quality generated sam-  
129 ples. Conditional WGAN (cWGAN) (Yu et al.,  
130 2019) extends WGAN by incorporating conditional  
131 variables, allowing the generation of samples con-  
132 ditioned on specific attributes and enhancing the  
133 model’s ability to generate more targeted and di-  
134 verse data samples. Vector Quantized Variational  
135 Autoencoders (VQ-VAE) (Van Den Oord et al.,  
136 2017) and VQ-GAN employ discrete latent repre-  
137 sentations through vector quantization. VQ-VAE  
138 improves its ability to handle complex data distri-  
139 butions compared to standard VAEs by learning a  
140 finite set of embeddings. VQ-GAN combines the  
141 VQ-VAE method with a discriminator to differenti-  
142 ate between real and fake data more effectively by  
143 learning a codebook.

144 In the realm of NLP, autoencoders are frequently  
145 combined to generate data in an embedding space  
146 (Malandrakis et al., 2019; Piedboeuf and Langlais,  
147 2022). This study leverages the VQ-GAN method  
148 to generate synthetic embeddings. Additionally,  
149 we analyze the semantic content of the synthetic  
150 embeddings produced by VQ-TEGAN and com-  
151 pare it with the embeddings created by mixup and  
152 the original text embedding data.

### 153 2.2 Text Augmentation

154 Text augmentation aims to improve model per-  
155 formance when data is insufficient. Early work  
156 includes EDA (Wei and Zou, 2019) and AEDA  
157 (Karimi et al., 2021). EDA employs four straight-  
158 forward data augmentation techniques: random  
159 swap, random insertion, random deletion, and syn-  
160 onym replacement. Similarly, AEDA operates by  
161 randomly inserting punctuation marks. TreeMix  
162 (Zhang et al., 2022) utilizes constituency parsing  
163 trees to decompose sentences into component sub-  
164 structures, which are then recombined using the  
165 mixup data augmentation method to generate new

166 sentences.

167 Instead of reorganizing words or sentences, an-  
168 other approach involves generating new text data  
169 using LLMs for data augmentation (Anaby-Tavor  
170 et al., 2020; Body et al., 2021; Dopierre et al.,  
171 2021; Ng et al., 2020). LAMBADA (Anaby-Tavor  
172 et al., 2020) fine-tunes a GPT model (Radford et al.,  
173 2019) on a small dataset and then augments it with  
174 the given label. BF-Translation (Body et al., 2021)  
175 uses the Google Translate API, with German as  
176 an intermediate language, to back-translate text for  
177 sentiment analysis. ProtAugment (Dopierre et al.,  
178 2021) combines paraphrases generated from the  
179 BART model with sentences produced through tradi-  
180 tional back-translation, improving intent detec-  
181 tion models via unsupervised meta-learning. This  
182 method utilizes paraphrasing-based data augmen-  
183 tation. SSMBA (Ng et al., 2020) is a word-level  
184 data augmentation technique that employs a corrup-  
185 tion function to mask specific tokens in a sentence  
186 and replace them with new tokens using a BERT  
187 model.

188 Furthermore, data augmentation in continuous  
189 embedding spaces, such as EmbedHalluc (Jian  
190 et al., 2022), has shown promising results. Specifi-  
191 cally, graph-based methods (Onan, 2023) and con-  
192 trastive learning (Jia et al., 2023) have been ex-  
193 plored for text augmentation. Embedding Aug-  
194 menter (Pellicer et al., 2023) is a technique that  
195 uses a word-changing algorithm with the GloVe  
196 model (Pennington et al., 2014) with 300 dimen-  
197 sions to find the most similar words.

198 This study investigates the use of synthetic em-  
199 beddings for data augmentation, where embeddings  
200 are derived from synonyms and related words. In  
201 particular, the proposed VQ-TEGAN model offers  
202 the advantage of being relatively lightweight com-  
203 pared to larger language models.

## 204 2.3 Fine-tuning of Pre-trained Language 205 Models

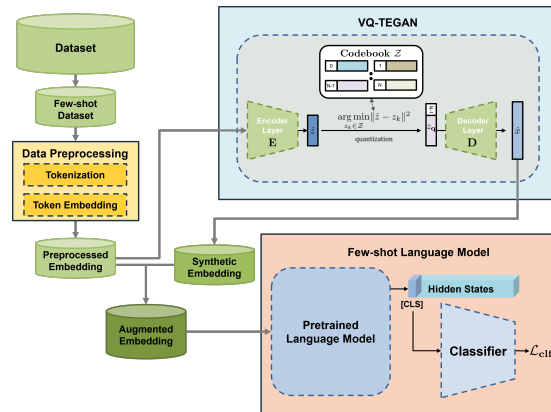
206 Numerous studies suggest using general models  
207 to address NLP problems (Kim, 2014; Huang  
208 et al., 2015; Kowsari et al., 2019). However, with  
209 the recent emergence of PLMs (*e.g.*, BERT and  
210 RoBERTa), there has been a surge in research on  
211 few-shot learning to leverage limited data with the  
212 help of PLMs (Gupta et al., 2020; Zhong et al.,  
213 2021; Chada and Natarajan, 2021; Ram et al.,  
214 2021). Some studies have applied data augmenta-  
215 tion to NLP classification tasks to improve few-shot  
216 learning performance (Wei et al., 2021; Jian et al.,

217 2022; Zhang et al., 2022; Jia et al., 2023). How-  
218 ever, the approach of creating new synthetic word  
219 embeddings for each word in a sentence, merging  
220 them, and using the resulting synthetic sentence  
221 embedding as training data for text classification  
222 has not yet been explored. In this context, we pro-  
223 pose VQ-TEGAN, the first attempt to apply the  
224 VQ-GAN method to generate new synthetic text  
225 embeddings for fine-tuning PLMs.

## 226 3 Methods

### 227 3.1 Overview

228 This research aims to evaluate the effectiveness  
229 of VQ-TEGAN in few-shot learning compared to  
230 benchmarks by performing classification tasks in  
231 limited data environments. The complete process  
232 of fine-tuning the PLM is illustrated in Figure 2.



233 Figure 2: Few-shot learning process using VQ-TEGAN

234 To preserve the integrity and diversity of the  
235 dataset, non-duplicating samples are randomly se-  
236 lected from each class for each classification task  
237 in the training and validation sets, respectively. The  
238 conversion of few-shot datasets to real embeddings  
239 is achieved using the PLM’s individual tokenizer  
240 and token embeddings, which are subsequently used  
241 to form preprocessed embeddings. The real em-  
242 beddings of the training set are then utilized to  
243 create synthetic embeddings through the pre-  
244 trained VQ-TEGAN. The synthetic embeddings  
245 for each real embedding are subsequently mixed to  
246 form the final augmented embeddings. The final  
247 augmented dataset, which includes one synthetic  
248 data point corresponding to each real data point,  
249 is used for few-shot learning. This approach takes  
250 advantage of the diversity introduced by the aug-  
251 mented data, operating under the assumption that  
252 it will enhance the learning capacity of the model  
when dealing with a restricted dataset (Arthaud

et al., 2021; Xie et al., 2020). Also, the freezing of word embeddings within the PLM during few-shot learning preserves the semantic integrity of the augmented dataset within the embeddings. This method proficiently transmits the intended semantics of the augmented dataset in few-shot learning contexts.

### 3.2 VQ-TEGAN

The architecture of a new generative model for text embedding data, VQ-TEGAN, is presented in detail in Figure 3. The primary objective is to train VQ-TEGAN directly within word embeddings in PLM to generate high-quality synthetic text embeddings. This approach has the advantage of leveraging PLM embeddings, eliminating the need for a separate training dataset. Furthermore, VQ-TEGAN allows the encapsulation of word embeddings with analogous attributes into quantized vectors, ensuring that the generated synthetic embeddings retain their distinct characteristics. The amount of training data depends on the number of word embeddings in PLMs. Note that RoBERTa-large and BERT-large have 50,265 and 30,522 embedding vectors, respectively. This approach has the advantage of utilizing embeddings of PLM, eliminating the need for a separate training dataset.

In VQ-VAE, a discrete-dimensional encoder output paired with an autoregressive decoder effectively solves the posterior collapse problem (Van Den Oord et al., 2017). VQ-TEGAN employs a similar structure to reconstruct the real embedding ( $x$ ) as the synthetic embedding ( $\hat{x}$ ) through the encoder  $\mathbf{E}$  - decoder  $\mathbf{D}$  structure illustrated in Figure 3. The input vector  $x \in \mathbb{R}^{n_x}$ , where  $n_x$  is the dimensionality of the input embedding, is compressed by the encoder  $\mathbf{E}$  into the latent vector  $\hat{z} \in \mathbb{R}^{n_z}$ , where  $n_z$  is the dimensionality of the codebook vector.

The latent vector  $\hat{z}$  is converted into one of the nearest codebook vectors,  $z_q \in \mathcal{Z}$ , by finding the distance to the vectors in the predefined discrete codebook, where  $\mathcal{Z} = \{z_k\}_{k=1}^K \subset \mathbb{R}^{n_z}$  and  $K$  is the number of codebook vectors. Specifically,  $\hat{z}$  is created from  $x$  and then quantized by replacing  $\hat{z}$  with the nearest codebook to obtain  $z_q$  such that:

$$z_q = \mathbf{q}(\hat{z}) := \arg \min_{z_k \in \mathcal{Z}} \|\hat{z} - z_k\|^2 \in \mathbb{R}^{n_z} \quad (1)$$

where  $\hat{z} = \mathbf{E}(x)$ . The reconstruction  $\hat{x} \approx x$  is given by:

$$\hat{x} = \mathbf{D}(z_q) \quad (2)$$

Backpropagation is not differentiable due to the quantization operation in Eq. 1. However, the model and codebook can be learned end-to-end via a loss function using a straight-through gradient estimator (Bengio et al., 2013) that copies the gradient from the decoder to the encoder as follows:

$$\mathcal{L}_{\text{VQ}}(\mathbf{E}, \mathbf{D}, \mathcal{Z}) = \|x - \hat{x}\| + 1 - \sigma(\hat{x}, x) + \|\text{sg}[\mathbf{E}(x)] - z_q\|^2 + \beta \times \|\text{sg}[z_q] - \mathbf{E}(x)\|^2 \quad (3)$$

Note that  $\|x - \hat{x}\|$  is a reconstruction loss ( $\mathcal{L}_{\text{rec}}$ );  $1 - \sigma(\hat{x}, x)$  is the cosine loss ( $\mathcal{L}_{\text{cos}}$ ) (Barz and Denzler, 2020) where  $\sigma(\cdot)$  represents the cosine similarity operation; and  $\|\text{sg}[z_q] - \mathbf{E}(x)\|^2$  is the commitment loss (Van Den Oord et al., 2017) where  $\text{sg}[\cdot]$  represents the stop-gradient operation.

To customize a learning approach for text embeddings, we modify the loss function commonly used in computer vision (Esser et al., 2021). Specifically, we replace the  $L_2$  loss with the  $L_1$  loss in  $\mathcal{L}_{\text{rec}}$ , a technique known for its effectiveness in high-resolution image restoration tasks (Zhao et al., 2016; Wu et al., 2021; Liu et al., 2021). The importance of cosine similarity in semantic analysis is derived from the inherent nature of text data embedding (Rahutomo et al., 2012; Pellicer et al., 2023).  $\mathcal{L}_{\text{cos}}$  is employed to ensure that the synthetic embedding  $\hat{x}$  is generated in a space characterized by high cosine similarity to the real embedding  $x$ .

The discriminator of VQ-TEGAN,  $\mathbf{Disc}$ , is responsible for distinguishing between real and fake embedding, resulting in a loss  $\mathcal{L}_{\text{Disc}}$  that follows the WGAN loss to efficiently train the generator (Arjovsky et al., 2017):

$$\mathcal{L}_{\text{GAN}}(\{\mathbf{E}, \mathbf{D}, \mathcal{Z}\}, \mathbf{Disc}) = \mathbf{Disc}(x) - \mathbf{Disc}(\hat{x}) \quad (4)$$

The complete objective to identify the optimal compression model  $\mathcal{Q}^* = \{\mathbf{E}^*, \mathbf{D}^*, \mathcal{Z}^*\}$  can be expressed as follows:

$$\mathcal{Q}^* = \arg \min_{\mathbf{E}, \mathbf{D}, \mathcal{Z}} \max_{\mathbf{Disc}} \mathbb{E}_{x \sim p(x)} [\mathcal{L}_{\text{VQ}}(\mathbf{E}, \mathbf{D}, \mathcal{Z}) + \mathcal{L}_{\text{GAN}}(\{\mathbf{E}, \mathbf{D}, \mathcal{Z}\}, \mathbf{Disc})] \quad (5)$$

VQ-TEGAN stands out for its scalability and memory efficiency in text embedding data augmentation, optimizing computational resources. The model’s parameters remain almost constant despite an increase in codebooks, growing only slightly from 5.03M (19.22MB) for 1024 codebooks to 5.42M (20.72MB) for 4096 codebooks. This

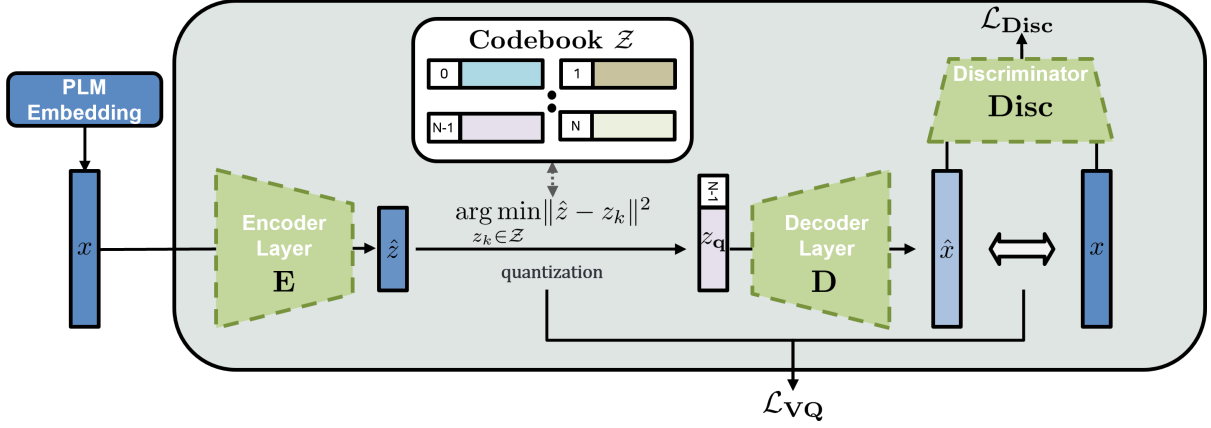


Figure 3: Model architectures of VQ-TEGAN

lightweight nature allows VQ-TEGAN to be deployed on various hardware, from high-end servers to resource-limited edge devices. Its compact design makes it ideal for scenarios that require robust text embedding augmentation without compromising performance. Training procedures are detailed in Appendix A.

### 3.3 Mixup Embedding

Mixup for word embedding, an application method devised by Guo et al. (2019), involves the linear interpolation of real and synthetic embeddings. We apply the mixup method as follows:

$$\tilde{x} = \lambda x + (1 - \lambda)\hat{x} \quad (6)$$

The mixup ratio  $\lambda$  specifies the proportion of real embedding ( $x$ ) in the mixed embedding. For instance, a  $\lambda$  of 1.0 indicates that the mixed embedding  $\tilde{x}$  is entirely composed of  $x$ , while a  $\lambda$  of 0.4 produces a mixture of 40% of  $x$  and 60% of  $\hat{x}$ . When  $\lambda$  is 0.0,  $\tilde{x}$  is composed of  $\hat{x}$  only.

## 4 Results & Discussions

### 4.1 Dataset

The research employs nine classification tasks from the GLUE benchmark dataset (Wang et al., 2018). The GLUE benchmark encompasses diverse tasks, including grammatical acceptability (CoLA), sentiment analysis (SST-2), sentence semantic equivalence (MRPC), semantic similarity (QQP), logical inference (MNLI-m, MNLI-mm), validity of sentence answers to questions (QNLI), and logical entailment (RTE), pronoun resolution (WNLI).

We randomly select 16 train and validation samples per class from the train and validation set of each task. The evaluations are based on the average results of five different seeds in the test set.

### 4.2 Impact of Cosine Loss

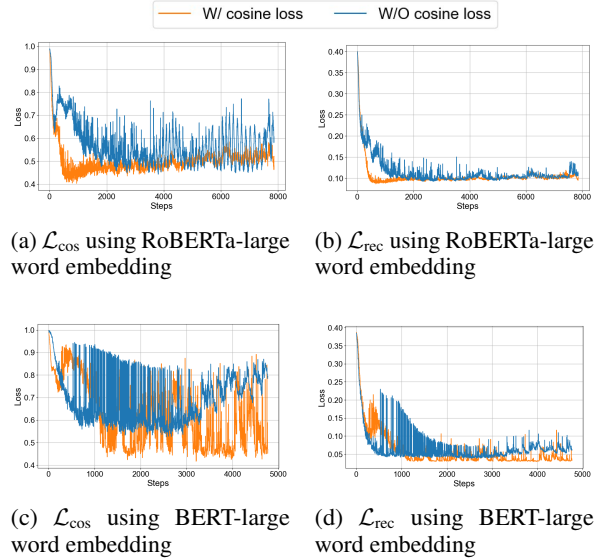


Figure 4:  $\mathcal{L}_{\cos}$  and  $\mathcal{L}_{\text{rec}}$  as a loss function with or without  $\mathcal{L}_{\cos}$  when training VQ-TEGAN

The analysis of Figure 4 underscores the importance of integrating the cosine loss term,  $\mathcal{L}_{\cos}$ , within Eq. 3. The integration stabilizes and accelerates the convergence, thus enhancing model performance in similarity measures and improving the quality of the reconstructed data.

Figures 4a and 4c illustrate the effect of cosine loss on the cosine similarity between the real embedding  $x$  and the synthetic embedding  $\hat{x}$ . The figures demonstrate that incorporating cosine loss in the generator’s loss function,  $\mathcal{L}_{VQ}$ , leads to faster and more stable convergence (orange line) compared to the method without cosine loss (blue line) during VQ-TEGAN training.

Figures 4b and 4d illustrate the reconstruction loss,  $\mathcal{L}_{\text{rec}}$ , defined as the  $L_1$  loss that quantifies the

362 difference between real and synthetic embeddings. 413  
363 Incorporation of cosine loss yields lower and more 414  
364 stable  $L_1$  loss values, indicating that synthetic em- 415  
365 beddings increasingly approximate the real input 416  
366 data. This observation implies that cosine loss en- 417  
367 hances the generator’s proficiency in accurately re- 418  
368 constructing inputs, thereby improving the overall 419  
369 fidelity of the generated embeddings. 420

370 Figure 4 shows that the word embeddings de- 421  
371 rived from RoBERTa-large demonstrate a more 422  
372 consistent convergence in comparison to those of 423  
373 BERT-large during the training phase. This obser- 424  
374 vation suggests that RoBERTa-large embeddings 425  
375 are more appropriate for training VQ-TEGAN, 426  
376 with the potential to produce embeddings that are 427  
377 semantically richer than those obtained from BERT- 428  
378 large embeddings. 429

### 379 4.3 Classification Performance in Few-shot 425 380 Learning 426

381 Table 1 provides a comprehensive analysis of the 427  
382 efficacy of various data augmentation methods, 428  
383 namely EDA, EmbedHalluc, and VQ-TEGAN, 429  
384 when implemented in few-shot learning scenar- 430  
385 ios using RoBERTa-large and BERT-large models 431  
386 across nine distinct tasks. The hyperparameters 432  
387 for few-shot learning are presented in Appendix B, 433  
388 while the benchmark methods are described in Ap- 434  
389 pendix C. The findings indicate that VQ-TEGAN 435  
390 consistently surpasses the other methods in most 436  
391 tasks, underscoring its robustness in text data aug- 437  
392 mentation. In particular, VQ-TEGAN significantly 438  
393 outperforms in seven tasks, with the exception 439  
394 of QNLI and RTE. However, VQ-TEGAN still 440  
395 achieves parity with EmbedHalluc on QNLI and is 441  
396 only 0.72% less accurate than EDA on RTE. 442

397 Although VQ-TEGAN demonstrates enhance- 443  
398 ments in RoBERTa-large, its performance remains 444  
399 comparable to other benchmarks when evaluated 445  
400 with BERT-large. EDA exhibits superior perfor- 446  
401 mance in MRPC (F1) with a score of 1.52, whereas 447  
402 EmbedHalluc surpasses in MNLI-mm, RTE, and 448  
403 WNLI by margins of 0.04%, 0.08%, and 0.94%, 449  
404 respectively. VQ-TEGAN also shows improved 450  
405 results, albeit marginally, with an increase of 0.08 451  
406 in CoLA (Matt.), and 0.38%, 0.02%, and 1.94% 452  
407 in SST-2, MNLI-m, and QNLI, respectively. It is 453  
408 important to note that no significant performance 454  
409 disparities are observed when these models are ap- 455  
410 plied to BERT-large. 456

411 In conclusion, VQ-TEGAN consistently sur- 457  
412 passes EDA and EmbedHalluc, particularly when 458

413 integrated with RoBERTa-large as opposed to 414  
415 BERT-large. The magnitude and complexity of 416  
417 the word embeddings of the employed PLM can 417  
418 significantly influence the extent of performance 418  
419 enhancement achieved with VQ-TEGAN. Given 419  
420 that VQ-TEGAN is trained directly on the word 420  
421 embeddings of the PLM, the utilization of more 421  
422 diverse and intricate embeddings for training cul- 422  
423 minates in more effective data augmentation. Con- 423  
424 sequently, VQ-TEGAN can be seen as a suitable 424  
425 data augmentation method to enhance the perfor- 425  
426 mance of larger PLMs relative to smaller ones. 426

### 427 4.4 Semantic Analysis on Mixup Embedding 425

427 Table 2 presents the three most prominent words 426  
428 decoded from the word embeddings of RoBERTa- 427  
429 large and BERT-large, demonstrating the highest 428  
430 cosine similarity to the mixup embeddings with 429  
431 different mixup ratio,  $\lambda$ . The words “beautiful”, 430  
432 “bad”, “characters”, and “doubts” are used as input, 431  
433 and the results illustrate the alterations in embed- 432  
434 dings under varying degrees of mixup. Note that 433  
435 the embeddings are congruent with the real embed- 434  
436 ding at  $\lambda = 1.0$ . The result is a representation of 435  
437 the words in the embeddings that demonstrate the 436  
438 highest cosine similarity to the real embedding for 437  
439 each PLM. The results show that the embeddings 438  
440 of all terms exhibit the highest cosine similarity to 439  
441 the embeddings of synonyms or capitalized forms 440  
442 for both PLMs. 441

442 In the case of the RoBERTa-large model, the list 442  
443 of closest word embeddings from  $\lambda = 0.8$  is iden- 443  
444 tical or slightly modified from  $\lambda = 1.0$ , including 444  
445 minor modifications to words (e.g., “suspicions”) 445  
446 or capitalization (e.g., “BAD”). That is, the seman- 446  
447 tic properties of the closest embeddings exhibit 447  
448 minimal variation relative to the case with  $\lambda = 1.0$ . 448  
449 When  $\lambda$  is set to 0.6, new words different from 449  
450 the list of  $\lambda = 1.0$  start to appear in the third rank 450  
451 (e.g., “magnificent” for the word “beautiful” and 451  
452 “lousy” for the word “bad”). As  $\lambda$  decreases to 0.4, 452  
453 many words that have similar semantic properties 453  
454 emerge in the list (e.g., “crappy” for the word “bad” 454  
455 and “protagonists” and “superheroes” for the word 455  
456 “characters”). This phenomenon is strengthened 456  
457 when  $\lambda$  decreases to 0.2, showing an increasing de- 457  
458 viation from the original words. For instance, the 458  
459 top three words decoded from the RoBERTa-large 459  
460 are “superheroes”, “mystic”, and “villan” for the 460  
461 word “characters”. For a  $\lambda$  of 0.0, the embedding is 461  
462 populated with novel words that are not related to 462  
463 the original words. It emphasizes the necessity of 463

Model	CoLA (Matt.)	SST-2 (acc)	MRPC (F1)	QQP (acc)	MNLI-m (acc)	MNLI-mm (acc)	QNLI (acc)	RTE (acc)	WNLI (acc)
<b>RoBERTa-large</b>	17.20±10.28	72.58±9.59	67.86±7.83	62.26±6.91	33.62±0.70	34.78±0.58	47.80±1.49	49.68±1.22	57.38±5.20
+EDA	12.42±6.78	70.48±6.78	68.68±13.98	57.22±18.77	33.78±1.22	33.88±2.34	49.22±1.24	<b>50.96</b> ±0.72	53.56±5.96
+EmbedHalluc	21.90±8.57	75.82±6.48	69.52±4.77	63.12±4.89	33.38±1.14	34.96±0.85	<b>49.64</b> ±0.75	49.54±1.01	55.32±8.20
+VQ-TEGAN	<b>29.66</b> ±8.02	<b>78.14</b> ±8.13	<b>72.50</b> ±4.16	<b>70.98</b> ±8.10	<b>34.68</b> ±1.14	<b>36.00</b> ±2.51	<b>49.64</b> ±1.40	50.24±0.42	<b>62.60</b> ±2.95
<b>BERT-large</b>	8.18±4.04	75.36±8.16	64.42±14.51	59.12±8.69	32.32±1.00	33.46±2.29	48.92±1.66	49.56±0.43	47.80±9.28
+EDA	10.48±3.59	78.22±4.36	<b>73.14</b> ±6.50	46.12±13.08	32.56±1.06	32.42±1.59	49.84±2.95	49.62±1.58	52.46±9.70
+EmbedHalluc	12.30±7.19	74.10±7.56	63.84±16.07	59.26±4.70	34.30±1.75	<b>35.12</b> ±2.21	48.60±2.30	<b>49.64</b> ±0.73	<b>53.68</b> ±8.11
+VQ-TEGAN	<b>12.38</b> ±4.53	<b>78.60</b> ±4.38	71.62±6.92	<b>66.98</b> ±5.59	<b>34.32</b> ±1.18	35.08±2.78	<b>51.78</b> ±1.27	49.56±0.50	52.74±6.77

Table 1: A comparative analysis of Conventional Fine-tuning, EDA, EmbedHalluc, and VQ-TEGAN, using RoBERTa-large and BERT-large as base models. The superior performance for each task is denoted in bold.

Word Embedding		RoBERTa-large				BERT-large			
$\lambda$	Rank	beautiful	bad	characters	doubts	beautiful	bad	characters	doubts
1.0	1	beautiful	bad	characters	doubts	beautiful	bad	characters	doubts
	2	gorgeous	Bad	character	doubt	gorgeous	good	character	doubted
	3	lovely	terrible	Characters	doubted	lovely	badly	protagonists	doubt
0.8	1	beautiful	bad	characters	doubts	beautiful	bad	characters	doubts
	2	gorgeous	Bad	character	doubted	gorgeous	badly	character	doubted
	3	lovely	BAD	Characters	suspicious	lovely	295	protagonists	doubt
0.6	1	beautiful	bad	characters	doubts	beautiful	bad	characters	doubts
	2	gorgeous	BAD	character	doubted	gorgeous	295	protagonists	[unused306]
	3	magnificent	lousy	Characters	suspicious	1738	321	1743	[unused298]
0.4	1	beautiful	bad	characters	doubts	1736	1736	1736	doubts
	2	gorgeous	lousy	protagonists	doubted	1732	276	1743	[unused659]
	3	magnificent	crappy	superheroes	suspicious	1738	326	1732	[unused276]
0.2	1	Beautiful	intertwined	superheroes	doubts	1736	1736	1736	[unused659]
	2	magnificent	sandy	mystic	timid	1732	276	1743	[unused80]
	3	the	crafted	vilains	dismay	1743	1732	1732	[unused176]
0.0	1	ACE	unfold	mystic	mystic	1736	1736	1736	[unused659]
	2	Apex	crafted	wretched	wretched	1732	1732	1732	[unused80]
	3	EA	intertwined	timid	timid	45th	45th	45th	[unused176]

Table 2: The top three words decoded from word embeddings in RoBERTa-large and BERT-large, exhibiting the highest degree of cosine similarity to the mixup embeddings with different  $\lambda$ .

the mixup for the augmentation via VQ-TEGAN.

In the case of the BERT-large model, at  $\lambda = 0.8$ , a minor change is observed for the word “bad”, but no change is observed for other words. Interestingly, the new words included in the word “bad” include the semantically unrelated word “295”. As  $\lambda$  decreases to 0.6, there is a significant increase in unrelated tokens and numbers observed, indicating a stronger deviation from the original words. As the value of  $\lambda$  is reduced from 0.4 to 0.0, the list is filled with semantically irrelevant words.

Our analysis indicates that as  $\lambda$  decreases, the mixup embeddings exhibit an increasing divergence from the original words. Furthermore, the mixup embeddings produced by RoBERTa-large are observed to encapsulate more semantically rich and contextually pertinent words at smaller  $\lambda$  compared to those generated by BERT-large. This observation suggests that the mixup embeddings of RoBERTa-large maintain a higher degree of semantic coherence under mixup conditions compared to BERT-large. This is corroborated by the classification performance presented in Table 1, which demonstrates that RoBERTa-large exhibits a sig-

nificant improvement in performance with mixup embeddings, whereas BERT-large does not show a comparable enhancement.

In conclusion, when VQ-TEGAN generates meaningful synthetic embeddings and integrates mixup embeddings with real embeddings for few-shot learning, it has the potential to facilitate the application of mixup embeddings with an expanded and more heterogeneous semantic spectrum for few-shot learning. Additional semantic analysis on mixup embeddings can be found in the Appendix D.

#### 4.5 Sensitivity Analysis on Mixup Ratio

In Table 3, we present a comparative analysis of the results derived from conventional fine-tuning and our proposed model, employing three distinct  $\lambda$  values (0.0, 0.2, and 0.4). The scenario with  $\lambda = 1$  was omitted from the sensitivity analysis due to its redundancy in merely duplicating the real embedding. Likewise, scenarios with  $\lambda = 0.6$  and  $\lambda = 0.8$  were excluded as their results did not show significant deviations from those presented in Table 2.

Model	CoLA (Matt.)	SST-2 (acc)	MRPC (F1)	QQP (acc)	MNLI-m (acc)	MNLI-mm (acc)	QNLI (acc)	RTE (acc)	WNLI (acc)
<b>RoBERTa-large</b>	17.20±10.28	72.58±9.59	67.86±7.83	62.26±6.91	33.62±0.70	34.78±0.58	47.80±1.49	49.68±1.22	57.38±5.20
w/ $\lambda = 0.0$	<b>28.32</b> ±11.60	<b>78.14</b> ±8.13	<b>71.98</b> ±6.36	<b>70.98</b> ±8.10	<b>34.60</b> ±1.56	<b>36.00</b> ±2.51	<b>48.22</b> ±1.42	<b>50.24</b> ±0.42	<b>60.96</b> ±4.42
w/ $\lambda = 0.2$	<b>29.66</b> ±8.02	<b>76.84</b> ±5.88	<b>72.50</b> ±4.16	<b>66.68</b> ±7.68	<b>34.40</b> ±0.90	34.66±1.07	<b>49.64</b> ±1.40	<b>50.02</b> ±0.65	<b>59.18</b> ±4.08
w/ $\lambda = 0.4$	<b>18.30</b> ±3.72	<b>74.42</b> ±7.33	<b>71.62</b> ±6.86	<b>65.24</b> ±9.24	<b>34.68</b> ±1.14	<b>35.60</b> ±2.29	<b>48.94</b> ±0.27	<b>50.22</b> ±0.76	<b>62.60</b> ±2.95
<b>BERT-large</b>	8.18±4.04	75.36±8.16	64.42±14.51	59.12±8.69	32.32±1.00	33.46±2.29	48.92±1.66	49.56±0.43	47.80±9.28
w/ $\lambda = 0.0$	<b>9.44</b> ±6.84	<b>77.34</b> ±5.00	<b>68.20</b> ±11.32	<b>66.98</b> ±5.59	<b>34.32</b> ±1.18	<b>35.08</b> ±2.78	<b>50.26</b> ±1.81	49.42±0.44	<b>52.74</b> ±6.77
w/ $\lambda = 0.2$	<b>12.38</b> ±4.53	<b>77.00</b> ±5.36	<b>71.62</b> ±6.92	<b>62.76</b> ±12.57	<b>33.46</b> ±1.57	<b>34.24</b> ±1.41	<b>50.12</b> ±0.98	49.46±0.91	<b>51.22</b> ±7.04
w/ $\lambda = 0.4$	<b>9.62</b> ±7.42	<b>78.60</b> ±4.38	<b>69.96</b> ±7.08	<b>63.78</b> ±7.14	<b>33.58</b> ±1.65	<b>34.02</b> ±3.10	<b>51.78</b> ±1.27	<b>49.56</b> ±0.50	<b>52.34</b> ±6.52

Table 3: A comparative analysis of conventional fine-tuning and VQ-TEGAN for different  $\lambda$ , using RoBERTa-large and BERT-large as base models. The bold numbers indicate instances where VQ-TEGAN outperforms conventional fine-tuning for each respective task, while underlined numbers indicate the highest performance.

Using RoBERTa-large for few-shot learning, VQ-TEGAN demonstrates superior performance relative to fine-tuning across all evaluated tasks. In general,  $\lambda = 0.0$  and  $\lambda = 0.2$  exhibit increased efficacy compared to traditional fine-tuning and  $\lambda = 0.4$ , with the exception of MNLI-m and WNLI. Specifically, for tasks such as SST-2, QQP, MNLI-mm, and RTE, the optimal results are observed with  $\lambda = 0.0$ . In contrast,  $\lambda = 0.2$  achieves superior results in CoLA, MRPC, and QNLI. In particular,  $\lambda = 0.4$  surpasses  $\lambda = 0.0$  and  $\lambda = 0.2$  exclusively in MNLI-m and WNLI. These findings indicate that the incorporation of synthetic embeddings or mixup embeddings significantly enhances model generalization and performance.

In contrast, using BERT-large for few-shot learning reveals a distinct pattern. Specifically, a  $\lambda$  value of 0.2 enhances performance beyond traditional fine-tuning in the CoLA and MRPC datasets. The most substantial performance improvements are achieved with  $\lambda = 0.4$  in the SST-2, QNLI, and RTE tasks. In particular, a  $\lambda$  value of 0.0 yields the highest performance metrics in QQP, MNLI-m, MNLI-mm, and WNLI. These observations suggest that the efficacy of BERT is differentially influenced by varying  $\lambda$  values and synthetic embeddings contingent on the specific task, thereby indicating the absence of a universally optimal  $\lambda$  value across all tasks.

## 5 Conclusion

This study introduces VQ-TEGAN, a novel data augmentation method for text embedding. VQ-TEGAN generates embeddings across various semantic and synonymic dimensions of PLM embeddings, facilitating more efficient and effective acquisition of a broader spectrum of semantics during the fine-tuning of PLMs with limited training datasets. Our empirical analysis reveals that

VQ-TEGAN (1) achieves superior performance enhancements on GLUE benchmark tasks in few-shot learning contexts, (2) is more compact and lightweight compared to other language models employed for data augmentation, (3) augments PLM performance, particularly when utilized with PLMs possessing larger embeddings, and (4) introduces a more efficient loss function for text embedding generation via the convergence of loss functions.

## 6 Limitations

Despite its novelty, there are limitations that need to be addressed in future work. As discussed in section 4.4, the semantic analysis of the closest PLM word embeddings to the mixup embeddings elucidates the potential for formulating a novel embedding space conducive to few-shot learning. However, a limitation is identified where VQ-TEGAN-generated embeddings may converge within a space similar to other semantic embeddings, attributable to the anisotropy issue inherent in PLM word embeddings (Ethayarajh, 2019; Li et al., 2020). A possible approach is to train VQ-TEGAN utilizing word embeddings derived from PLMs that have been refined through contrastive learning (Gao et al., 2021), addressing the anisotropy issue within the embedding space. Lastly, this study exclusively investigates the instances of VQ-TEGAN utilizing RoBERTa-large and BERT-large. For subsequent study, a broader spectrum of PLMs should be explored for the implementation of VQ-TEGAN.

## 7 Ethics Statement

This paper investigates data augmentation in the generation of embeddings for few-shot learning. It is not anticipated that this research will raise any ethical or social issues. All data utilized in this study is publicly accessible and has been utilized by numerous researchers. The proposed method



586	does not introduce any ethical biases into the data.	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. <i>arXiv preprint arXiv:1810.04805</i> .	637 638 639 640
587	<b>References</b>		
588	Ateret Anaby-Tavor, Boaz Carmeli, Esther Goldbraich, Amir Kantor, George Kour, Segev Shlomov, Naama Tepper, and Naama Zwerdling. 2020. Do not have enough data? deep learning to the rescue! In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pages 7383–7390.	Thomas Dopierre, Christophe Gravier, and Wilfried Logerais. 2021. Protaugment: Intent detection meta-learning through unsupervised diverse paraphrasing. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 2454–2466. Association for Computational Linguistics.	641 642 643 644 645 646 647 648
589			
590			
591			
592			
593			
594	Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein generative adversarial networks. In <i>International conference on machine learning</i> , pages 214–223. PMLR.	Zach Eaton-Rosen, Felix J. S. Bragman, Sébastien Ourselin, and M. Jorge Cardoso. 2018. Improving data augmentation for medical image segmentation. In <i>OpenReview</i> .	649 650 651 652
595			
596			
597			
598	Farid Arthaud, Rachel Bawden, and Alexandra Birch. 2021. Few-shot learning through contextual data augmentation. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 1049–1062, Online. Association for Computational Linguistics.	Yousef El-Laham and Svitlana Vyetrenko. 2022. Style-time: Style transfer for synthetic time series generation. In <i>Proceedings of the Third ACM International Conference on AI in Finance</i> , pages 489–496.	653 654 655 656
599			
600			
601			
602			
603			
604	Bjorn Barz and Joachim Denzler. 2020. Deep learning on small datasets without pre-training using cosine loss. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pages 1371–1380.	Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming transformers for high-resolution image synthesis. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 12873–12883.	657 658 659 660 661
605			
606			
607			
608			
609	Yoshua Bengio, Nicholas Léonard, and Aaron Courville. 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. <i>arXiv preprint arXiv:1308.3432</i> .	Kawin Ethayarajh. 2019. How contextual are contextualized word representations? comparing the geometry of bert, elmo, and gpt-2 embeddings. <i>arXiv preprint arXiv:1909.00512</i> .	662 663 664 665
610			
611			
612			
613	Thomas Body, Xiaohui Tao, Yuefeng Li, Lin Li, and Ning Zhong. 2021. Using back-and-forth translation to create artificial augmented textual data for sentiment analysis models. <i>Expert Systems with Applications</i> , 178:115033.	Elizabeth Fons, Paula Dawson, Xiao-jun Zeng, John Keane, and Alexandros Iosifidis. 2020. Evaluating data augmentation for financial time series classification. <i>arXiv preprint arXiv:2010.15111</i> .	666 667 668 669
614			
615			
616			
617			
618	Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. <i>Advances in neural information processing systems</i> , 33:1877–1901.	Tianyu Gao, Adam Fisch, and Danqi Chen. 2020. Making pre-trained language models better few-shot learners. <i>arXiv preprint arXiv:2012.15723</i> .	670 671 672
619			
620			
621			
622			
623			
624	Rakesh Chada and Pradeep Natarajan. 2021. Fewshotqa: A simple framework for few-shot learning of question answering tasks using pre-trained text-to-text models. <i>arXiv preprint arXiv:2109.01951</i> .	Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. <i>arXiv preprint arXiv:2104.08821</i> .	673 674 675
625			
626			
627			
628	Jiaao Chen, Zichao Yang, and Diyi Yang. 2020a. Mix-text: Linguistically-informed interpolation of hidden space for semi-supervised text classification. <i>arXiv preprint arXiv:2004.12239</i> .	Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. <i>Advances in neural information processing systems</i> , 27.	676 677 678 679 680
629			
630			
631			
632	Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020b. A simple framework for contrastive learning of visual representations. In <i>International conference on machine learning</i> , pages 1597–1607. PMLR.	Hongyu Guo, Yongyi Mao, and Richong Zhang. 2019. Augmenting data with mixup for sentence classification: An empirical study. <i>arXiv preprint arXiv:1905.08941</i> .	681 682 683 684
633			
634			
635			
636			
637		Aakriti Gupta, Kapil Thadani, and Neil O’Hare. 2020. Effective few-shot classification with transfer learning. In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 1061–1066.	685 686 687 688 689

690	Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional lstm-crf models for sequence tagging. <i>arXiv preprint arXiv:1508.01991</i> .	742
691		743
692		744
693	Ouyang Jia, Huimin Huang, Jiaxin Ren, Luodi Xie, and Yinyin Xiao. 2023. Contrastive learning with text augmentation for text classification. <i>Applied Intelligence</i> , pages 1–10.	745
694		746
695		747
696		748
697	Yiren Jian, Chongyang Gao, and Soroush Vosoughi. 2022. <a href="#">Embedding hallucination for few-shot language fine-tuning</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , 3, pages 5522–5530. Association for Computational Linguistics.	749
698		750
699		751
700		752
701		753
702		754
703		755
704		756
705	Akbar Karimi, Leonardo Rossi, and Andrea Prati. 2021. <a href="#">AEDA: An easier data augmentation technique for text classification</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 2748–2754, Punta Cana, Dominican Republic. Association for Computational Linguistics.	757
706		758
707		759
708		760
709		
710	Justin Ker, Lipo Wang, Jai Rao, and Tchoyoson Lim. 2017. Deep learning applications in medical image analysis. <i>Ieee Access</i> , 6:9375–9389.	761
711		762
712		763
713	Yoon Kim. 2014. Convolutional neural networks for sentence classification. <i>arXiv preprint arXiv:1408.5882</i> .	764
714		765
715		766
716	Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. <i>arXiv preprint arXiv:1312.6114</i> .	767
717		768
718		769
719	Kamran Kowsari, Kiana Jafari Meimandi, Mojtaba Heidarysafa, Sanjana Mendu, Laura Barnes, and Donald Brown. 2019. Text classification algorithms: A survey. <i>Information</i> , 10(4):150.	770
720		771
721		772
722		773
723	Bohan Li, Hao Zhou, Junxian He, Mingxuan Wang, Yiming Yang, and Lei Li. 2020. On the sentence embeddings from pre-trained language models. <i>arXiv preprint arXiv:2011.05864</i> .	774
724		775
725		776
726		777
727	Yang Liu, Zhenyue Qin, Saeed Anwar, Pan Ji, Dongwoo Kim, Sabrina Caldwell, and Tom Gedeon. 2021. Invertible denoising network: A light solution for real noise removal. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 13365–13374.	778
728		779
729		780
730		781
731		782
732		783
733	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	784
734		785
735		786
736		787
737		788
738	Shayne Longpre, Yu Wang, and Christopher DuBois. 2020. How effective is task-agnostic data augmentation for pretrained transformers? <i>arXiv preprint arXiv:2010.01764</i> .	789
739		790
740		791
741		792
		793
		794
	Nikolaos Malandrakis, Minmin Shen, Anuj Goyal, Shuyang Gao, Abhishek Sethi, and Angeliki Metallinou. 2019. <a href="#">Controlled text generation for data augmentation in intelligent artificial agents</a> . In <i>Proceedings of the 3rd Workshop on Neural Generation and Translation</i> , pages 90–98, Hong Kong. Association for Computational Linguistics.	795
		796
	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. <i>arXiv preprint arXiv:1301.3781</i> .	797
		798
	Nathan Ng, Kyunghyun Cho, and Marzyeh Ghassemi. 2020. Ssmba: Self-supervised manifold based data augmentation for improving out-of-domain robustness. <i>arXiv preprint arXiv:2009.10195</i> .	799
		800
	Aytuğ Onan. 2023. Gtr-ga: Harnessing the power of graph-based neural networks and genetic algorithms for text augmentation. <i>Expert Systems with Applications</i> , page 120908.	801
		802
	Lucas Francisco Amaral Orosco Pellicer, Taynan Maier Ferreira, and Anna Helena Reali Costa. 2023. Data augmentation techniques in natural language processing. <i>Applied Soft Computing</i> , 132:109803.	803
		804
	Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In <i>Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)</i> , pages 1532–1543.	805
		806
	Frédéric Piedboeuf and Philippe Langlais. 2022. Effective data augmentation for sentence classification using one vae per class. In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 3454–3464.	807
		808
	Lutz Prechelt. 2002. Early stopping-but when? In <i>Neural Networks: Tricks of the trade</i> , pages 55–69. Springer.	809
		810
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	811
		812
	Faisal Rahutomo, Teruaki Kitasuka, and Masayoshi Arimitsu. 2012. Semantic cosine similarity. In <i>The 7th international student conference on advanced science and technology ICAST</i> , volume 4, page 1.	813
		814
	Ori Ram, Yuval Kirstain, Jonathan Berant, Amir Globerson, and Omer Levy. 2021. Few-shot question answering by pretraining span selection. <i>arXiv preprint arXiv:2101.00438</i> .	815
		816
	Marc’Aurelio Ranzato, Fu Jie Huang, Y-Lan Boureau, and Yann LeCun. 2007. <a href="#">Unsupervised learning of invariant feature hierarchies with applications to object recognition</a> . In <i>2007 IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 1–8.	817
		818
		819
		820

795 Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural  
796 discrete representation learning. *Advances in neural*  
797 *information processing systems*, 30.

798 Alex Wang, Amanpreet Singh, Julian Michael, Felix  
799 Hill, Omer Levy, and Samuel R Bowman. 2018.  
800 Glue: A multi-task benchmark and analysis platform  
801 for natural language understanding. *arXiv preprint*  
802 *arXiv:1804.07461*.

803 William Yang Wang and Diyi Yang. 2015. That’s so an-  
804 noying!!!: A lexical and frame-semantic embedding  
805 based data augmentation approach to automatic cat-  
806 egorization of annoying behaviors using# petpeeve  
807 tweets. In *Proceedings of the 2015 conference on*  
808 *empirical methods in natural language processing*,  
809 pages 2557–2563.

810 Jason Wei, Chengyu Huang, Soroush Vosoughi,  
811 Yu Cheng, and Shiqi Xu. 2021. Few-shot text  
812 classification with triplet networks, data augmen-  
813 tation, and curriculum learning. *arXiv preprint*  
814 *arXiv:2103.07552*.

815 Jason Wei and Kai Zou. 2019. Eda: Easy data augmen-  
816 tation techniques for boosting performance on text clas-  
817 sification tasks. *arXiv preprint arXiv:1901.11196*.

818 Haiyan Wu, Yanyun Qu, Shaohui Lin, Jian Zhou, Ruizhi  
819 Qiao, Zhizhong Zhang, Yuan Xie, and Lizhuang Ma.  
820 2021. Contrastive learning for compact single image  
821 dehazing. In *Proceedings of the IEEE/CVF Confer-*  
822 *ence on Computer Vision and Pattern Recognition*,  
823 pages 10551–10560.

824 Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and  
825 Quoc Le. 2020. Unsupervised data augmentation for  
826 consistency training. *Advances in neural information*  
827 *processing systems*, 33:6256–6268.

828 Ying Yu, Bingying Tang, Ronglai Lin, Shufa Han, Tang  
829 Tang, and Ming Chen. 2019. *Cwgan: Conditional*  
830 *wasserstein generative adversarial nets for fault data*  
831 *generation*. In *2019 IEEE International Conference*  
832 *on Robotics and Biomimetics (ROBIO)*, pages 2713–  
833 2718.

834 Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and  
835 David Lopez-Paz. 2017. mixup: Beyond empirical  
836 risk minimization. *arXiv preprint arXiv:1710.09412*.

837 Le Zhang, Zichao Yang, and Diyi Yang. 2022. Treemix:  
838 Compositional constituency-based data augmentation  
839 for natural language understanding. *arXiv preprint*  
840 *arXiv:2205.06153*.

841 Hang Zhao, Orazio Gallo, Iuri Frosio, and Jan Kautz.  
842 2016. Loss functions for image restoration with neu-  
843 ral networks. *IEEE Transactions on computational*  
844 *imaging*, 3(1):47–57.

845 Wanjun Zhong, Duyu Tang, Jiahai Wang, Jian Yin, and  
846 Nan Duan. 2021. Useradapter: Few-shot user learn-  
847 ing in sentiment analysis. In *Findings of the Associ-*  
848 *ation for Computational Linguistics: ACL-IJCNLP*  
849 *2021*, pages 1484–1488.

## A Training Details for VQ-TEGAN

850 The generator architecture includes an encoder, a  
851 decoder, and a codebook of latent vectors. The en-  
852 coder is composed of four sequential blocks, each  
853 containing a fully connected layer, batch normaliza-  
854 tion, and the LeakyReLU activation function (Jian  
855 et al., 2022). This encoder progressively reduces  
856 the dimensionality to 1024, 512, 256, and 128. The  
857 codebook comprises quantized latent vectors that  
858 correspond to the output dimensions of the encoder.  
859 The quantity of codebook vectors is adjusted as  
860 a hyperparameter during the experimental proce-  
861 dures. The decoder, which structurally parallels  
862 the encoder, consists of four blocks that expand the  
863 quantized codebook vectors to dimensions of 128,  
864 256, 512, and 1024. The discriminator is structured  
865 with three blocks, having dimensions of 512, 512,  
866 and 1, respectively, and produces a singular tensor  
867 output. VQ-TEGAN is subjected to training for 10  
868 epochs with a batch size of 64, utilizing the Adam  
869 optimizer ( $\beta = (0.5, 0.999)$ ) and a fixed random  
870 seed of 42. The training process includes a grid  
871 search for the learning rates of  $2e^{-5}$  and  $5e^{-5}$ , as  
872 well as codebook vector quantities of 1024, 2048,  
873 and 4096. 874

## B Hyperparameters for Few-shot Learning

875 The model is trained using learning rates of  $1e^{-5}$   
876 and  $2e^{-5}$ , with batch sizes of 4 and 8. Random  
877 number generation seeds of 13, 21, 42, 87, and 100  
878 are utilized. The training process was capped at 150  
879 epochs, with the final model being selected based  
880 on validation accuracy at each epoch. An early  
881 stopping mechanism is used to mitigate overfitting,  
882 halting training if no improvement in validation  
883 accuracy is observed after 100 epochs (Prechelt,  
884 2002). 885

886 To train the PLM with augmented embeddings,  
887 comprehensive experiments are conducted across  
888 all parameters. The mixup ratios for  $x$  and  $\hat{x}$  are  
889 evaluated at  $\lambda$  values of 0, 0.2, and 0.4 as illus-  
890 trated in Eq. 6. Both EDA and EmbedHalluc are  
891 executed using default settings, with EDA’s data  
892 augmentation further explored by generating 4 and  
893 9 additional samples. 894

895 The algorithms are implemented using Python  
896 3.10.8 and PyTorch 1.13.1. The experiments are  
897 carried out on an Ubuntu 20.04.6 system equipped  
898 with a Nvidia RTX 3090 TI (24 GB RAM) and an  
899 Intel(R) Xeon(R) Gold 6226R CPU @ 2.90GHz.

The NLTK 3.8.1 toolkit is used for synonym replacement in the EDA process. RoBERTa-large and BERT-large models, along with their tokenizers, are sourced from the Hugging Face Transformers library.

## C Benchmarks

The performance of VQ-TEGAN is evaluated in comparison to established benchmarks: conventional fine-tuning, EDA, and EmbedHalluc based on cWGAN.

- **Conventional Fine-tuning** constitutes a fundamental approach where a few-shot language model is trained exclusively on the provided dataset, devoid of any supplementary data augmentation.
- **EDA**(Wei and Zou, 2019) represents a data augmentation that incorporates four principal techniques: synonym replacement, random deletion, random swap, and random addition. This method is both intuitive and efficient, facilitating the generation of a substantial number of synthetic sentences in a straightforward manner.
- **EmbedHalluc**(Jian et al., 2022) leverages cWGAN to augment textual data within the embedding space. The training process encompasses both the generator and the discriminator, with data augmentation being realized by doubling the few-shot data via the generator.

## D Text Embedding Analysis

Tables 4, 5, 6, and 7 present the words of the embedding tokens that exhibit the three highest cosine similarities between the mixup embeddings generated by VQ-TEGAN and the embeddings within the PLM under various parameter configurations. Tables 4 and 5 elucidate the results for RoBERTa-large embeddings, while Tables 6 and 7 illustrate the results for BERT-large embeddings. The parameter configurations encompass the learning rate of VQ-TEGAN, the number of codebook vectors, and five distinct values of  $\lambda$ .

Tables 4 and 5 show the evolving patterns in the semantic representations of the embeddings as the parameter  $\lambda$  changes, as discerned through our comprehensive analysis.

At a  $\lambda$  value of 0.0, where the embeddings are synthesized exclusively by VQ-TEGAN in the absence of any real embeddings, the cosine similarity fails to effectively discern relatedness or syn-

onymy. This observation implies that the synthetic embeddings may exhibit abstract or non-traditional associations at this  $\lambda$ , which deviate from the conventional semantic relationships observed in real embeddings.

A notable change is observed when the parameter  $\lambda$  is elevated to 0.2 or 0.4. At these  $\lambda$  values, the top three synonyms for each text sample exhibit increased diversity, which means that the mixed embeddings now encapsulate a wider spectrum of semantic similarities. For instance, Table 4 demonstrates that the mixup embedding of “beautiful”, with a  $\lambda$  of 0.2, achieves the highest cosine similarities of 0.751, 0.742, and 0.740 with the embeddings of adjectives bearing analogous meanings, such as “exquisite”, “magnificent”, and “marvellous”, respectively, for  $lr_{VQ}$  of 5e-05 and a codebook size of 4096. In Table 5, it is evident that the mixup embedding of “characters” manifests the three highest cosine similarities of 0.670, 0.670, and 0.658 with the embeddings of nouns possessing similar or identical meanings, such as “villains”, “superheroes”, and “characters”. This observation is pivotal, as it suggests that the mixup embedding at these  $\lambda$  values transcends a mere replication of the original meanings. Moreover, it introduces an extensive array of related concepts with real embeddings. The inclusion of 20% of  $x$  functions as an anchor, anchoring the synthetic embedding within the original semantic framework while still allowing the introduction of novel nuances. This equilibrium, facilitated by synthetic embeddings generated by VQ-TEGAN, which adeptly constructs the semantic space of PLM’s embedding, culminates in mixup embeddings that are enriched with supplementary contextual meaning. This generates a more nuanced and comprehensive semantic comprehension.

Upon observation, it was observed that when the parameter  $\lambda$  exceeds the threshold of 0.6, the augmented embeddings exhibit a pronounced resemblance to the real embedding  $x$ . This phenomenon indicates that at elevated values, the mixup embeddings converge more closely with the semantic attributes of the genuine embedding, thereby diminishing the distinctions from the original embedding. As a result, this convergence may precipitate issues of data redundancy, as the mixup embeddings may not provide substantially novel or diverse information relative to the original dataset.

Tables 6 and 7 demonstrate that the evaluation of BERT-large mixup embeddings via cosine sim-

Table 4: The text and cosine similarity metrics of the three words decoded from RoBERTa-large, exhibiting the highest cosine similarity between the word embeddings in RoBERTa-large and the mixup embeddings using the adjectives ‘beautiful’ and ‘bad’. The parameter  $\text{lr}_{VQ}$  denotes the learning rate of the VQ-TEGAN, while the term codebook refers to the number of codebook vectors in VQ-TEGAN ( $K$  of  $\mathcal{Z} = \{z_k\}_{k=1}^K$  as specified in Eq. 1). Additionally,  $\lambda$  represents the rate of authentic text embedding.

Parameter	Codebook	Text $\lambda$	beautiful			bad			
			Rank 1	Rank 2	Rank 3	Rank 1	Rank 2	Rank 3	
2e-05	1024	0.0	ACE <sub>(0.691)</sub>	Apex <sub>(0.688)</sub>	EA <sub>(0.687)</sub>	unfold <sub>(0.606)</sub>	crafted <sub>(0.602)</sub>	intertwined <sub>(0.601)</sub>	
		0.2	Beautiful <sub>(0.702)</sub>	magnificent <sub>(0.701)</sub>	the <sub>(0.697)</sub>	intertwined <sub>(0.621)</sub>	sandy <sub>(0.619)</sub>	crafted <sub>(0.618)</sub>	
		0.4	beautiful <sub>(0.83)</sub>	gorgeous <sub>(0.769)</sub>	magnificent <sub>(0.764)</sub>	bad <sub>(0.732)</sub>	lousy <sub>(0.660)</sub>	crappy <sub>(0.658)</sub>	
		0.6	beautiful <sub>(0.935)</sub>	gorgeous <sub>(0.834)</sub>	lovely <sub>(0.790)</sub>	bad <sub>(0.882)</sub>	BAD <sub>(0.683)</sub>	lousy <sub>(0.682)</sub>	
		0.8	beautiful <sub>(0.987)</sub>	gorgeous <sub>(0.854)</sub>	lovely <sub>(0.804)</sub>	bad <sub>(0.974)</sub>	Bad <sub>(0.713)</sub>	BAD <sub>(0.707)</sub>	
	2048	0.0	extravagant <sub>(0.721)</sub>	stereotypical <sub>(0.718)</sub>	astounding <sub>(0.718)</sub>	extravagant <sub>(0.721)</sub>	stereotypical <sub>(0.718)</sub>	astounding <sub>(0.718)</sub>	
		0.2	exquisite <sub>(0.751)</sub>	magnificent <sub>(0.742)</sub>	astounding <sub>(0.740)</sub>	obnoxious <sub>(0.728)</sub>	extravagant <sub>(0.723)</sub>	stereotypical <sub>(0.723)</sub>	
		0.4	beautiful <sub>(0.819)</sub>	magnificent <sub>(0.787)</sub>	gorgeous <sub>(0.781)</sub>	bad <sub>(0.815)</sub>	crappy <sub>(0.731)</sub>	horrendous <sub>(0.731)</sub>	
		0.6	beautiful <sub>(0.923)</sub>	gorgeous <sub>(0.840)</sub>	magnificent <sub>(0.799)</sub>	bad <sub>(0.922)</sub>	BAD <sub>(0.728)</sub>	terrible <sub>(0.719)</sub>	
		0.8	beautiful <sub>(0.983)</sub>	gorgeous <sub>(0.860)</sub>	lovely <sub>(0.806)</sub>	bad <sub>(0.983)</sub>	BAD <sub>(0.721)</sub>	Bad <sub>(0.719)</sub>	
4096	1024	0.0	Shoes <sub>(0.702)</sub>	Changes <sub>(0.702)</sub>	Moments <sub>(0.695)</sub>	hateful <sub>(0.725)</sub>	despicable <sub>(0.722)</sub>	wretched <sub>(0.718)</sub>	
		0.2	Shoes <sub>(0.703)</sub>	Hair <sub>(0.700)</sub>	Awesome <sub>(0.699)</sub>	horrendous <sub>(0.749)</sub>	hateful <sub>(0.741)</sub>	despicable <sub>(0.740)</sub>	
		0.4	Beautiful <sub>(0.758)</sub>	beautiful <sub>(0.743)</sub>	gorgeous <sub>(0.707)</sub>	bad <sub>(0.866)</sub>	BAD <sub>(0.749)</sub>	horrendous <sub>(0.746)</sub>	
		0.6	beautiful <sub>(0.891)</sub>	gorgeous <sub>(0.807)</sub>	Beautiful <sub>(0.779)</sub>	bad <sub>(0.950)</sub>	BAD <sub>(0.749)</sub>	terrible <sub>(0.734)</sub>	
		0.8	beautiful <sub>(0.976)</sub>	gorgeous <sub>(0.852)</sub>	lovely <sub>(0.797)</sub>	bad <sub>(0.99)</sub>	BAD <sub>(0.725)</sub>	terrible <sub>(0.719)</sub>	
	5e-05	1024	0.0	demonic <sub>(0.637)</sub>	anarchist <sub>(0.630)</sub>	superhero <sub>(0.626)</sub>	demonic <sub>(0.637)</sub>	anarchist <sub>(0.630)</sub>	superhero <sub>(0.626)</sub>
			0.2	marvelous <sub>(0.683)</sub>	majestic <sub>(0.680)</sub>	magnificent <sub>(0.680)</sub>	demonic <sub>(0.674)</sub>	hateful <sub>(0.668)</sub>	heinous <sub>(0.668)</sub>
			0.4	beautiful <sub>(0.822)</sub>	gorgeous <sub>(0.763)</sub>	magnificent <sub>(0.748)</sub>	bad <sub>(0.811)</sub>	BAD <sub>(0.700)</sub>	horrendous <sub>(0.691)</sub>
			0.6	beautiful <sub>(0.932)</sub>	gorgeous <sub>(0.832)</sub>	lovely <sub>(0.788)</sub>	bad <sub>(0.929)</sub>	BAD <sub>(0.730)</sub>	terrible <sub>(0.715)</sub>
			0.8	beautiful <sub>(0.986)</sub>	gorgeous <sub>(0.854)</sub>	lovely <sub>(0.804)</sub>	bad <sub>(0.986)</sub>	BAD <sub>(0.720)</sub>	terrible <sub>(0.714)</sub>
2048	2048	0.0	exquisite <sub>(0.647)</sub>	extravagant <sub>(0.636)</sub>	ludicrous <sub>(0.633)</sub>	exquisite <sub>(0.647)</sub>	extravagant <sub>(0.636)</sub>	ludicrous <sub>(0.633)</sub>	
		0.2	exquisite <sub>(0.718)</sub>	extravagant <sub>(0.680)</sub>	astounding <sub>(0.679)</sub>	troublesome <sub>(0.675)</sub>	ludicrous <sub>(0.671)</sub>	exquisite <sub>(0.670)</sub>	
		0.4	beautiful <sub>(0.789)</sub>	exquisite <sub>(0.762)</sub>	magnificent <sub>(0.750)</sub>	bad <sub>(0.771)</sub>	horrendous <sub>(0.706)</sub>	dreadful <sub>(0.703)</sub>	
		0.6	beautiful <sub>(0.915)</sub>	gorgeous <sub>(0.817)</sub>	magnificent <sub>(0.781)</sub>	bad <sub>(0.909)</sub>	terrible <sub>(0.712)</sub>	horrendous <sub>(0.703)</sub>	
		0.8	beautiful <sub>(0.982)</sub>	gorgeous <sub>(0.851)</sub>	lovely <sub>(0.798)</sub>	bad <sub>(0.981)</sub>	terrible <sub>(0.717)</sub>	horrible <sub>(0.703)</sub>	
	4096	0.0	esoteric <sub>(0.705)</sub>	clandestine <sub>(0.698)</sub>	subversive <sub>(0.697)</sub>	esoteric <sub>(0.705)</sub>	clandestine <sub>(0.698)</sub>	subversive <sub>(0.697)</sub>	
		0.2	exquisite <sub>(0.739)</sub>	magnificent <sub>(0.724)</sub>	marvelous <sub>(0.723)</sub>	nefarious <sub>(0.720)</sub>	obnoxious <sub>(0.713)</sub>	crappy <sub>(0.712)</sub>	
		0.4	beautiful <sub>(0.800)</sub>	magnificent <sub>(0.774)</sub>	exquisite <sub>(0.769)</sub>	bad <sub>(0.788)</sub>	crappy <sub>(0.736)</sub>	lousy <sub>(0.734)</sub>	
		0.6	beautiful <sub>(0.911)</sub>	gorgeous <sub>(0.829)</sub>	magnificent <sub>(0.794)</sub>	bad <sub>(0.907)</sub>	BAD <sub>(0.746)</sub>	crappy <sub>(0.726)</sub>	
		0.8	beautiful <sub>(0.979)</sub>	gorgeous <sub>(0.857)</sub>	lovely <sub>(0.807)</sub>	bad <sub>(0.979)</sub>	BAD <sub>(0.733)</sub>	Bad <sub>(0.724)</sub>	
1.0	beautiful <sub>(1.000)</sub>	gorgeous <sub>(0.846)</sub>	lovely <sub>(0.792)</sub>	bad <sub>(1.000)</sub>	Bad <sub>(0.706)</sub>	terrible <sub>(0.691)</sub>			

Table 5: The text and cosine similarity metrics of the three words decoded from RoBERTa-large, exhibiting the highest cosine similarity between the word embeddings in RoBERTa-large and the mixup embeddings using the nouns ‘characters’ and ‘doubts’

Parameter	Text	characters			doubts			
		Codebook	$\lambda$	Rank 1	Rank 2	Rank 3	Rank 1	Rank 2
2e-05	1024	0.0	mystic <sub>(0.719)</sub>	wretched <sub>(0.718)</sub>	timid <sub>(0.715)</sub>	mystic <sub>(0.719)</sub>	wretched <sub>(0.718)</sub>	timid <sub>(0.715)</sub>
	0.2	superheroes <sub>(0.719)</sub>	mystic <sub>(0.716)</sub>	villains <sub>(0.716)</sub>	doubts <sub>(0.752)</sub>	timid <sub>(0.725)</sub>	dismay <sub>(0.722)</sub>	
	0.4	characters <sub>(0.826)</sub>	protagonists <sub>(0.743)</sub>	superheroes <sub>(0.740)</sub>	doubts <sub>(0.877)</sub>	doubted <sub>(0.742)</sub>	suspicious <sub>(0.736)</sub>	
	0.6	characters <sub>(0.934)</sub>	character <sub>(0.744)</sub>	Characters <sub>(0.744)</sub>	doubts <sub>(0.954)</sub>	doubted <sub>(0.751)</sub>	suspicious <sub>(0.745)</sub>	
	0.8	characters <sub>(0.987)</sub>	character <sub>(0.764)</sub>	Characters <sub>(0.738)</sub>	doubts <sub>(0.991)</sub>	doubted <sub>(0.736)</sub>	suspicious <sub>(0.730)</sub>	
2048	0.0	extravagant <sub>(0.721)</sub>	stereotypical <sub>(0.718)</sub>	astounding <sub>(0.718)</sub>	extravagant <sub>(0.721)</sub>	stereotypical <sub>(0.718)</sub>	astounding <sub>(0.718)</sub>	
	0.2	stereotypical <sub>(0.724)</sub>	extravagant <sub>(0.719)</sub>	imaginative <sub>(0.713)</sub>	doubtful <sub>(0.722)</sub>	doubts <sub>(0.718)</sub>	bogus <sub>(0.718)</sub>	
	0.4	characters <sub>(0.787)</sub>	protagonists <sub>(0.729)</sub>	protagonist <sub>(0.715)</sub>	doubts <sub>(0.852)</sub>	doubted <sub>(0.752)</sub>	doubtful <sub>(0.743)</sub>	
	0.6	characters <sub>(0.915)</sub>	character <sub>(0.743)</sub>	protagonists <sub>(0.738)</sub>	doubts <sub>(0.942)</sub>	doubted <sub>(0.762)</sub>	suspicious <sub>(0.737)</sub>	
	0.8	characters <sub>(0.982)</sub>	character <sub>(0.767)</sub>	Characters <sub>(0.730)</sub>	doubts <sub>(0.988)</sub>	doubted <sub>(0.744)</sub>	suspicious <sub>(0.729)</sub>	
4096	0.0	hateful <sub>(0.725)</sub>	despicable <sub>(0.722)</sub>	wretched <sub>(0.718)</sub>	hateful <sub>(0.725)</sub>	despicable <sub>(0.722)</sub>	wretched <sub>(0.718)</sub>	
	0.2	demonic <sub>(0.718)</sub>	nonsensical <sub>(0.716)</sub>	despicable <sub>(0.715)</sub>	doubts <sub>(0.752)</sub>	doubtful <sub>(0.739)</sub>	bogus <sub>(0.724)</sub>	
	0.4	characters <sub>(0.839)</sub>	protagonists <sub>(0.733)</sub>	protagonist <sub>(0.724)</sub>	doubts <sub>(0.887)</sub>	doubted <sub>(0.761)</sub>	doubtful <sub>(0.748)</sub>	
	0.6	characters <sub>(0.944)</sub>	character <sub>(0.755)</sub>	protagonists <sub>(0.729)</sub>	doubts <sub>(0.961)</sub>	doubted <sub>(0.760)</sub>	suspicious <sub>(0.738)</sub>	
	0.8	characters <sub>(0.989)</sub>	character <sub>(0.767)</sub>	Characters <sub>(0.727)</sub>	doubts <sub>(0.993)</sub>	doubted <sub>(0.738)</sub>	suspicious <sub>(0.725)</sub>	
5e-05	1024	0.0	boxes <sub>(0.624)</sub>	cats <sub>(0.618)</sub>	demonic <sub>(0.637)</sub>	anarchist <sub>(0.630)</sub>	superhero <sub>(0.626)</sub>	
	0.2	Characters <sub>(0.703)</sub>	Characters <sub>(0.699)</sub>	protagonists <sub>(0.676)</sub>	doubts <sub>(0.690)</sub>	demonic <sub>(0.663)</sub>	baffled <sub>(0.658)</sub>	
	0.4	characters <sub>(0.874)</sub>	Characters <sub>(0.762)</sub>	protagonists <sub>(0.720)</sub>	doubts <sub>(0.853)</sub>	doubted <sub>(0.712)</sub>	suspicious <sub>(0.703)</sub>	
	0.6	characters <sub>(0.960)</sub>	Characters <sub>(0.764)</sub>	character <sub>(0.747)</sub>	doubts <sub>(0.949)</sub>	doubted <sub>(0.738)</sub>	suspicious <sub>(0.729)</sub>	
	0.8	characters <sub>(0.993)</sub>	character <sub>(0.760)</sub>	Characters <sub>(0.741)</sub>	doubts <sub>(0.990)</sub>	doubted <sub>(0.730)</sub>	suspicious <sub>(0.724)</sub>	
2048	0.0	exquisite <sub>(0.647)</sub>	extravagant <sub>(0.636)</sub>	ludicrous <sub>(0.633)</sub>	exquisite <sub>(0.647)</sub>	extravagant <sub>(0.636)</sub>	ludicrous <sub>(0.633)</sub>	
	0.2	exquisite <sub>(0.664)</sub>	extravagant <sub>(0.656)</sub>	outlandish <sub>(0.654)</sub>	ludicrous <sub>(0.669)</sub>	bogus <sub>(0.667)</sub>	exquisite <sub>(0.666)</sub>	
	0.4	characters <sub>(0.766)</sub>	protagonists <sub>(0.676)</sub>	protagonist <sub>(0.663)</sub>	doubts <sub>(0.814)</sub>	doubtful <sub>(0.709)</sub>	doubted <sub>(0.699)</sub>	
	0.6	characters <sub>(0.911)</sub>	Characters <sub>(0.726)</sub>	protagonists <sub>(0.708)</sub>	doubts <sub>(0.932)</sub>	doubted <sub>(0.737)</sub>	doubtful <sub>(0.714)</sub>	
	0.8	characters <sub>(0.982)</sub>	character <sub>(0.760)</sub>	Characters <sub>(0.704)</sub>	doubts <sub>(0.987)</sub>	doubted <sub>(0.733)</sub>	suspicious <sub>(0.719)</sub>	
4096	0.0	impooverished <sub>(0.634)</sub>	obsolete <sub>(0.628)</sub>	obscured <sub>(0.627)</sub>	impooverished <sub>(0.634)</sub>	obsolete <sub>(0.628)</sub>	obscured <sub>(0.627)</sub>	
	0.2	villains <sub>(0.670)</sub>	superheroes <sub>(0.670)</sub>	characters <sub>(0.658)</sub>	doubts <sub>(0.727)</sub>	disagreements <sub>(0.674)</sub>	shortcomings <sub>(0.674)</sub>	
	0.4	characters <sub>(0.823)</sub>	Characters <sub>(0.726)</sub>	protagonists <sub>(0.707)</sub>	doubts <sub>(0.866)</sub>	doubted <sub>(0.716)</sub>	suspicious <sub>(0.712)</sub>	
	0.6	characters <sub>(0.933)</sub>	Characters <sub>(0.753)</sub>	character <sub>(0.726)</sub>	doubts <sub>(0.951)</sub>	doubted <sub>(0.737)</sub>	suspicious <sub>(0.733)</sub>	
	0.8	characters <sub>(0.987)</sub>	character <sub>(0.755)</sub>	Characters <sub>(0.742)</sub>	doubts <sub>(0.990)</sub>	doubted <sub>(0.730)</sub>	suspicious <sub>(0.725)</sub>	
1.0	characters <sub>(1.000)</sub>	character <sub>(0.756)</sub>	Characters <sub>(0.711)</sub>	doubts <sub>(1.000)</sub>	doubt <sub>(0.718)</sub>	doubted <sub>(0.708)</sub>		

Table 6: The text and cosine similarity metrics of the three words decoded from BERT-large, exhibiting the highest cosine similarity between the word embeddings in RoBERTa-large and the mixup embeddings using the adjectives ‘beautiful’ and ‘bad’.

Parameter	Codebook	Text $\lambda$	beautiful			bad		
			Rank 1	Rank 2	Rank 3	Rank 1	Rank 2	Rank 3
2e-05	1024	0.0	1736 <sub>(0.642)</sub>	1732 <sub>(0.630)</sub>	45th <sub>(0.629)</sub>	1736 <sub>(0.642)</sub>	1732 <sub>(0.630)</sub>	45th <sub>(0.629)</sub>
		0.2	1736 <sub>(0.653)</sub>	1732 <sub>(0.641)</sub>	1743 <sub>(0.638)</sub>	1736 <sub>(0.651)</sub>	276 <sub>(0.640)</sub>	1732 <sub>(0.638)</sub>
		0.4	1736 <sub>(0.650)</sub>	1732 <sub>(0.639)</sub>	1738 <sub>(0.637)</sub>	1736 <sub>(0.635)</sub>	276 <sub>(0.633)</sub>	326 <sub>(0.632)</sub>
		0.6	beautiful <sub>(0.771)</sub>	gorgeous <sub>(0.616)</sub>	1738 <sub>(0.608)</sub>	bad <sub>(0.783)</sub>	295 <sub>(0.582)</sub>	321 <sub>(0.0.582)</sub>
		0.8	beautiful <sub>(0.938)</sub>	gorgeous <sub>(0.656)</sub>	lovely <sub>(0.596)</sub>	bad <sub>(0.948)</sub>	badly <sub>(0.467)</sub>	295 <sub>(0.463)</sub>
		0.0	1726 <sub>(0.328)</sub>	##た <sub>(0.319)</sub>	1777 <sub>(0.309)</sub>	1726 <sub>(0.328)</sub>	##た <sub>(0.319)</sub>	1777 <sub>(0.309)</sub>
		0.2	1726 <sub>(0.354)</sub>	##た <sub>(0.339)</sub>	1777 <sub>(0.338)</sub>	1726 <sub>(0.347)</sub>	1666 <sub>(0.330)</sub>	##た <sub>(0.330)</sub>
		0.4	1726 <sub>(0.386)</sub>	1777 <sub>(0.374)</sub>	1666 <sub>(0.373)</sub>	bad <sub>(0.383)</sub>	283 <sub>(0.377)</sub>	1744 <sub>(0.355)</sub>
		0.6	beautiful <sub>(0.606)</sub>	beautifully <sub>(0.448)</sub>	gorgeous <sub>(0.445)</sub>	bad <sub>(0.648)</sub>	283 <sub>(0.377)</sub>	326 <sub>(0.374)</sub>
		0.8	beautiful <sub>(0.884)</sub>	gorgeous <sub>(0.588)</sub>	lovely <sub>(0.536)</sub>	bad <sub>(0.908)</sub>	badly <sub>(0.418)</sub>	good <sub>(0.400)</sub>
4096	1024	0.0	1760 <sub>(0.483)</sub>	1709 <sub>(0.472)</sub>	1761 <sub>(0.472)</sub>	1760 <sub>(0.483)</sub>	1709 <sub>(0.472)</sub>	1761 <sub>(0.472)</sub>
		0.2	1760 <sub>(0.518)</sub>	44th <sub>(0.507)</sub>	##た <sub>(0.502)</sub>	271 <sub>(0.503)</sub>	1760 <sub>(0.503)</sub>	1709 <sub>(0.503)</sub>
		0.4	beautiful <sub>(0.605)</sub>	1760 <sub>(0.535)</sub>	1738 <sub>(0.528)</sub>	bad <sub>(0.612)</sub>	295 <sub>(0.518)</sub>	271 <sub>(0.514)</sub>
		0.6	beautiful <sub>(0.825)</sub>	gorgeous <sub>(0.594)</sub>	magnificent <sub>(0.551)</sub>	bad <sub>(0.845)</sub>	295 <sub>(0.480)</sub>	321 <sub>(0.469)</sub>
		0.8	beautiful <sub>(0.963)</sub>	gorgeous <sub>(0.635)</sub>	lovely <sub>(0.579)</sub>	bad <sub>(0.970)</sub>	good <sub>(0.442)</sub>	badly <sub>(0.432)</sub>
		0.0	## (0.771)	## (0.758)	1738 <sub>(0.754)</sub>	## (0.771)	## (0.758)	1738 <sub>(0.754)</sub>
		0.2	## (0.759)	## (0.754)	1738 <sub>(0.753)</sub>	## (0.743)	1738 <sub>(0.742)</sub>	1732 <sub>(0.738)</sub>
		0.4	beautiful <sub>(0.776)</sub>	1738 <sub>(0.695)</sub>	## (0.694)	bad <sub>(0.768)</sub>	295 <sub>(0.661)</sub>	283 <sub>(0.654)</sub>
		0.6	beautiful <sub>(0.917)</sub>	gorgeous <sub>(0.687)</sub>	magnificent <sub>(0.623)</sub>	bad <sub>(0.923)</sub>	295 <sub>(0.536)</sub>	283 <sub>(0.533)</sub>
		0.8	beautiful <sub>(0.984)</sub>	gorgeous <sub>(0.663)</sub>	lovely <sub>(0.597)</sub>	bad <sub>(0.986)</sub>	badly <sub>(0.448)</sub>	good <sub>(0.446)</sub>
2048	1024	0.0	1729 <sub>(0.485)</sub>	1744 <sub>(0.484)</sub>	1756 <sub>(0.481)</sub>	1729 <sub>(0.485)</sub>	1744 <sub>(0.484)</sub>	1756 <sub>(0.481)</sub>
		0.2	1729 <sub>(0.513)</sub>	1744 <sub>(0.510)</sub>	1734 <sub>(0.510)</sub>	1729 <sub>(0.512)</sub>	298 <sub>(0.505)</sub>	1732 <sub>(0.505)</sub>
		0.4	beautiful <sub>(0.610)</sub>	1738 <sub>(0.531)</sub>	1734 <sub>(0.530)</sub>	bad <sub>(0.609)</sub>	298 <sub>(0.516)</sub>	1729 <sub>(0.513)</sub>
		0.6	beautiful <sub>(0.819)</sub>	gorgeous <sub>(0.610)</sub>	lovely <sub>(0.550)</sub>	bad <sub>(0.836)</sub>	295 <sub>(0.479)</sub>	298 <sub>(0.476)</sub>
		0.8	beautiful <sub>(0.959)</sub>	gorgeous <sub>(0.645)</sub>	lovely <sub>(0.580)</sub>	bad <sub>(0.967)</sub>	good <sub>(0.441)</sub>	badly <sub>(0.433)</sub>
		0.0	1655 <sub>(0.417)</sub>	1717 <sub>(0.413)</sub>	##之 <sub>(0.412)</sub>	1655 <sub>(0.417)</sub>	1717 <sub>(0.413)</sub>	##之 <sub>(0.412)</sub>
		0.2	1655 <sub>(0.448)</sub>	29th <sub>(0.447)</sub>	1717 <sub>(0.442)</sub>	266 <sub>(0.444)</sub>	495 <sub>(0.439)</sub>	207 <sub>(0.437)</sub>
		0.4	beautiful <sub>(0.556)</sub>	1738 <sub>(0.475)</sub>	29th <sub>(0.474)</sub>	bad <sub>(0.567)</sub>	266 <sub>(0.462)</sub>	495 <sub>(0.462)</sub>
		0.6	beautiful <sub>(0.787)</sub>	gorgeous <sub>(0.562)</sub>	lovely <sub>(0.514)</sub>	bad <sub>(0.811)</sub>	283 <sub>(0.450)</sub>	304 <sub>(0.444)</sub>
		0.8	beautiful <sub>(0.952)</sub>	gorgeous <sub>(0.626)</sub>	lovely <sub>(0.568)</sub>	bad <sub>(0.961)</sub>	good <sub>(0.429)</sub>	badly <sub>(0.418)</sub>
1.0	beautiful <sub>(1.000)</sub>	gorgeous <sub>(0.619)</sub>	lovely <sub>(0.557)</sub>	bad <sub>(1.000)</sub>	good <sub>(0.450)</sub>	badly <sub>(0.401)</sub>		

Table 7: The text and cosine similarity metrics of the three words decoded from RoBERTa-large, exhibiting the highest cosine similarity between the word embeddings in RoBERTa-large and the mixup embeddings using the nouns ‘characters’ and ‘doubts’

Parameter	Text	characters					doubts				
		Codebook	$\lambda$	Rank 1	Rank 2	Rank 3	Rank 1	Rank 2	Rank 3	Rank 1	Rank 2
2e-05	1024	0.0	1736 <sub>(0.642)</sub>	1732 <sub>(0.630)</sub>	45th <sub>(0.629)</sub>	[unused659] <sub>(0.595)</sub>	[unused80] <sub>(0.593)</sub>	[unused176] <sub>(0.592)</sub>			
		0.2	1736 <sub>(0.650)</sub>	1743 <sub>(0.641)</sub>	1732 <sub>(0.640)</sub>	[unused659] <sub>(0.636)</sub>	[unused80] <sub>(0.633)</sub>	[unused176] <sub>(0.632)</sub>			
		0.4	1736 <sub>(0.639)</sub>	1743 <sub>(0.636)</sub>	1732 <sub>(0.632)</sub>	doubts <sub>(0.703)</sub>	[unused659] <sub>(0.662)</sub>	[unused276] <sub>(0.662)</sub>			
		0.6	characters <sub>(0.812)</sub>	protagonists <sub>(0.605)</sub>	1743 <sub>(0.591)</sub>	doubts <sub>(0.863)</sub>	[unused306] <sub>(0.658)</sub>	[unused298] <sub>(0.658)</sub>			
		0.8	characters <sub>(0.952)</sub>	character <sub>(0.622)</sub>	protagonists <sub>(0.581)</sub>	doubts <sub>(0.968)</sub>	doubt <sub>(0.679)</sub>	doubt <sub>(0.615)</sub>			
		0.0	1726 <sub>(0.328)</sub>	##た <sub>(0.319)</sub>	1777 <sub>(0.309)</sub>	1726 <sub>(0.328)</sub>	##た <sub>(0.319)</sub>	1777 <sub>(0.309)</sub>			
		0.2	1726 <sub>(0.358)</sub>	##た <sub>(0.338)</sub>	1744 <sub>(0.335)</sub>	1726 <sub>(0.378)</sub>	##た <sub>(0.366)</sub>	1666 <sub>(0.360)</sub>			
		0.4	characters <sub>(0.472)</sub>	1726 <sub>(0.393)</sub>	protagonist <sub>(0.384)</sub>	doubts <sub>(0.525)</sub>	1726 <sub>(0.436)</sub>	1679 <sub>(0.423)</sub>			
		0.6	characters <sub>(0.696)</sub>	protagonists <sub>(0.470)</sub>	protagonist <sub>(0.456)</sub>	doubts <sub>(0.748)</sub>	doubt <sub>(0.522)</sub>	[unused144] <sub>(0.496)</sub>			
		0.8	characters <sub>(0.915)</sub>	character <sub>(0.595)</sub>	protagonists <sub>(0.527)</sub>	doubts <sub>(0.936)</sub>	doubt <sub>(0.631)</sub>	doubt <sub>(0.587)</sub>			
4096	1024	0.0	1760 <sub>(0.483)</sub>	1709 <sub>(0.472)</sub>	1761 <sub>(0.472)</sub>	[unused26] <sub>(0.547)</sub>	[unused757] <sub>(0.546)</sub>				
		0.2	1760 <sub>(0.519)</sub>	44th <sub>(0.508)</sub>	1764 <sub>(0.507)</sub>	[unused26] <sub>(0.598)</sub>	[unused248] <sub>(0.596)</sub>				
		0.4	characters <sub>(0.626)</sub>	1760 <sub>(0.528)</sub>	1764 <sub>(0.527)</sub>	doubts <sub>(0.727)</sub>	[unused857] <sub>(0.630)</sub>				
		0.6	characters <sub>(0.848)</sub>	character <sub>(0.579)</sub>	protagonists <sub>(0.563)</sub>	doubts <sub>(0.883)</sub>	[unused298] <sub>(0.628)</sub>				
		0.8	characters <sub>(0.970)</sub>	character <sub>(0.642)</sub>	protagonists <sub>(0.550)</sub>	doubts <sub>(0.975)</sub>	doubt <sub>(0.618)</sub>				
		0.0	[unused467] <sub>(0.840)</sub>	[unused499] <sub>(0.836)</sub>	[unused962] <sub>(0.836)</sub>	[unused467] <sub>(0.840)</sub>	[unused962] <sub>(0.836)</sub>				
		0.2	[unused467] <sub>(0.807)</sub>	[unused257] <sub>(0.805)</sub>	[unused962] <sub>(0.804)</sub>	[unused306] <sub>(0.826)</sub>	[unused467] <sub>(0.826)</sub>				
		0.4	characters <sub>(0.791)</sub>	[unused257] <sub>(0.708)</sub>	[unused467] <sub>(0.708)</sub>	doubts <sub>(0.876)</sub>	[unused306] <sub>(0.770)</sub>				
		0.6	characters <sub>(0.925)</sub>	protagonists <sub>(0.638)</sub>	protagonist <sub>(0.596)</sub>	doubts <sub>(0.957)</sub>	[unused313] <sub>(0.768)</sub>				
		0.8	characters <sub>(0.986)</sub>	character <sub>(0.621)</sub>	Characters <sub>(0.572)</sub>	doubts <sub>(0.992)</sub>	[unused306] <sub>(0.687)</sub>				
2048	1024	0.0	1729 <sub>(0.485)</sub>	1744 <sub>(0.484)</sub>	1756 <sub>(0.481)</sub>	[unused158] <sub>(0.573)</sub>	[unused306] <sub>(0.572)</sub>				
		0.2	1744 <sub>(0.517)</sub>	1729 <sub>(0.516)</sub>	1734 <sub>(0.515)</sub>	[unused306] <sub>(0.626)</sub>	[unused439] <sub>(0.622)</sub>				
		0.4	characters <sub>(0.633)</sub>	1734 <sub>(0.535)</sub>	1744 <sub>(0.529)</sub>	doubts <sub>(0.748)</sub>	[unused2] <sub>(0.650)</sub>				
		0.6	characters <sub>(0.842)</sub>	character <sub>(0.569)</sub>	protagonists <sub>(0.559)</sub>	doubts <sub>(0.896)</sub>	[unused306] <sub>(0.645)</sub>				
		0.8	characters <sub>(0.967)</sub>	character <sub>(0.637)</sub>	protagonists <sub>(0.550)</sub>	doubts <sub>(0.978)</sub>	doubt <sub>(0.614)</sub>				
		0.0	1655 <sub>(0.417)</sub>	1717 <sub>(0.413)</sub>	##た <sub>(0.412)</sub>	[unused892] <sub>(0.578)</sub>	[unused949] <sub>(0.574)</sub>				
		0.2	1655 <sub>(0.456)</sub>	##た <sub>(0.452)</sub>	266 <sub>(0.449)</sub>	[unused892] <sub>(0.628)</sub>	[unused949] <sub>(0.627)</sub>				
		0.4	characters <sub>(0.577)</sub>	##成 <sub>(0.483)</sub>	##志 <sub>(0.707)</sub>	doubts <sub>(0.743)</sub>	[unused943] <sub>(0.654)</sub>				
		0.6	characters <sub>(0.813)</sub>	protagonists <sub>(0.526)</sub>	character <sub>(0.510)</sub>	doubts <sub>(0.895)</sub>	[unused298] <sub>(0.642)</sub>				
		0.8	characters <sub>(0.961)</sub>	character <sub>(0.614)</sub>	protagonists <sub>(0.539)</sub>	doubts <sub>(0.978)</sub>	doubt <sub>(0.616)</sub>				
5e-05	1024	0.0	[unused467] <sub>(0.840)</sub>	[unused499] <sub>(0.836)</sub>	[unused962] <sub>(0.836)</sub>	[unused467] <sub>(0.840)</sub>	[unused962] <sub>(0.836)</sub>				
		0.2	[unused467] <sub>(0.807)</sub>	[unused257] <sub>(0.805)</sub>	[unused962] <sub>(0.804)</sub>	[unused306] <sub>(0.826)</sub>	[unused467] <sub>(0.826)</sub>				
		0.4	characters <sub>(0.791)</sub>	[unused257] <sub>(0.708)</sub>	[unused467] <sub>(0.708)</sub>	doubts <sub>(0.876)</sub>	[unused306] <sub>(0.770)</sub>				
		0.6	characters <sub>(0.925)</sub>	protagonists <sub>(0.638)</sub>	protagonist <sub>(0.596)</sub>	doubts <sub>(0.957)</sub>	[unused306] <sub>(0.687)</sub>				
		0.8	characters <sub>(0.986)</sub>	character <sub>(0.621)</sub>	Characters <sub>(0.572)</sub>	doubts <sub>(0.992)</sub>	doubt <sub>(0.631)</sub>				
		0.0	1729 <sub>(0.485)</sub>	1744 <sub>(0.484)</sub>	1756 <sub>(0.481)</sub>	[unused158] <sub>(0.573)</sub>	[unused306] <sub>(0.572)</sub>				
		0.2	1744 <sub>(0.517)</sub>	1729 <sub>(0.516)</sub>	1734 <sub>(0.515)</sub>	[unused306] <sub>(0.626)</sub>	[unused439] <sub>(0.622)</sub>				
		0.4	characters <sub>(0.633)</sub>	1734 <sub>(0.535)</sub>	1744 <sub>(0.529)</sub>	doubts <sub>(0.748)</sub>	[unused2] <sub>(0.650)</sub>				
		0.6	characters <sub>(0.842)</sub>	character <sub>(0.569)</sub>	protagonists <sub>(0.559)</sub>	doubts <sub>(0.896)</sub>	[unused306] <sub>(0.645)</sub>				
		0.8	characters <sub>(0.967)</sub>	character <sub>(0.637)</sub>	protagonists <sub>(0.550)</sub>	doubts <sub>(0.978)</sub>	doubt <sub>(0.614)</sub>				
4096	1024	0.0	1655 <sub>(0.417)</sub>	1717 <sub>(0.413)</sub>	##た <sub>(0.412)</sub>	[unused892] <sub>(0.578)</sub>	[unused943] <sub>(0.574)</sub>				
		0.2	1655 <sub>(0.456)</sub>	##た <sub>(0.452)</sub>	266 <sub>(0.449)</sub>	[unused892] <sub>(0.628)</sub>	[unused943] <sub>(0.627)</sub>				
		0.4	characters <sub>(0.577)</sub>	##成 <sub>(0.483)</sub>	##志 <sub>(0.707)</sub>	doubts <sub>(0.743)</sub>	[unused943] <sub>(0.654)</sub>				
		0.6	characters <sub>(0.813)</sub>	protagonists <sub>(0.526)</sub>	character <sub>(0.510)</sub>	doubts <sub>(0.895)</sub>	[unused298] <sub>(0.642)</sub>				
		0.8	characters <sub>(0.961)</sub>	character <sub>(0.614)</sub>	protagonists <sub>(0.539)</sub>	doubts <sub>(0.978)</sub>	doubt <sub>(0.616)</sub>				
		1.0	characters <sub>(1.000)</sub>	character <sub>(0.647)</sub>	protagonists <sub>(0.502)</sub>	doubts <sub>(1.000)</sub>	doubt <sub>(0.634)</sub>				



1001 ilarity indicates that those embeddings with the  
1002 highest cosine similarity exhibit inferior semantic  
1003 coherence compared to RoBERTa-large.

1004 When  $\lambda$  is set to 0, the mixup embeddings  
1005 occupy a space distinct from the real embed-  
1006 dings. This phenomenon arises due to the non-  
1007 convergence of the cosine similarity depicted in  
1008 Figure 4c, despite the partial convergence of the  
1009 reconstruction loss illustrated in Figure 4d. Further-  
1010 more, it is apparent that the mixup embeddings are  
1011 predominantly characterized by synthetic embed-  
1012 dings when  $\lambda$  is 0.2 and 0.4, with only embeddings  
1013 being identified in a space similar to synthetic em-  
1014 beddings. For  $\lambda$  values of 0.6 and 0.8, the mixup  
1015 embeddings exhibit a greater resemblance to the  
1016 real embeddings, with a minority of embeddings  
1017 situated in a space similar to that of the synthetic  
1018 embeddings. This observation substantiates that  
1019 the mixup embeddings for BERT-large do not pos-  
1020 sess a semantic meaning that is markedly distinct  
1021 from the real and synthetic embeddings.

1022 Tables 4 and 5 illustrate that VQ-TEGAN gener-  
1023 ates semantic embeddings that provide RoBERTa-  
1024 large with access to a more diverse and meaningful  
1025 embedding space for learning. Conversely, Tables  
1026 6 and 7 reveal that the mixup embeddings on BERT-  
1027 large exhibit less significant cosine similarity com-  
1028 pared to those augmented on RoBERTa-large em-  
1029 beddings.